Characterizing License Practices in Maven Central and Their Relationship with CVEs: A CS846 Course Report

4/3/25

Haonan Zhang, Christina Li, Paul Wooseok Lee



Motivation & Background

- Open-source software is central to modern development.
- Maven Central hosts hundreds of thousands of artifacts.
- Licenses determine usage, reuse, and redistribution.
- Vulnerability (CVE) reporting is essential for security risk management.



What is the Relationship between Licenses and CVE Patterns?

Hypothesis

 Certain license types may influence collaboration and patching processes, potentially correlating with higher or lower vulnerability rates.

Objective

 Systematically analyze license data and CVEs in Maven Central to uncover patterns that can guide more informed artifact and license adoption choices.

Bradles on evils.
Understanding the Threats of Upstream Vulnerabilities to Downstream Projects in the
Maven Ecosystem Yuhan We ¹¹¹ , Zeliang Yu ¹¹¹ , Ming Wen ¹¹² , Qiang L ¹¹ , Deging Zeu ¹¹ , Hai Jin ²¹ ¹ Shool of Coher Science and Evolutioning, Mendebarg University of Science and Technology, Clina

Current Studies Primarily Focus on Either License or CVEs

Mitigating Persistence of Open-Source Vulnerabilities in Maven Ecosystem

Analyzing the Direct and Transitive Impact of Vulnerabilities onto Different Artifact Repositories

CS 846





CS 846

Hugwei China

Daniel A. Almeida¹ ⁽²⁾ · Gail C. Murphy¹ Greg Wilson² · Michael Hoye³

Investigating whether and how software developers

n the Adoption of Open Source Software Licensing - A Patter

Software Institute (Switzerland controlls amb writige

Collection

nderstand open source software licensing

PAGE 4



Research Questions

RQ1:

• What are the characteristics and trends of license adoption and CVE incidence across Maven Central artifacts?

RQ2:

• Do specific license types correlate with higher or lower vulnerability incidence in Maven Central artifacts?

CS 846

PAGE 5



WATERLOO

Dataset and Data Extraction



- Data sourced from Maven Central, Libraries.io, Neo4j, and Weaver API.
- Approximately 278,984 records extracted from a subset (over 3 days) out of 658K records.
- Data split into batches and processed via Python scripts using JSON formats.

PAGE 6

Tools and Methodology







WATERLOO

Integrating CVE Data

- CVE information retrieved via the Weaver API with cypher queries.
- Queries parameterized by library names allow integration of vulnerability data with library metadata.
- Data processed and aggregated using Python (json, collections.Counter).

Answering RQ1:

RQ1:

• What are the characteristics and trends of license adoption and CVE



Answering RQ1 Cont'd



Finding 1:

Artifacts under permissive licenses like the Apache-2.0 exhibit a **notably high** number of CVEs.
In contrast, licenses with stronger copyleft provisions, which have lower adoption rates, report **fewer** vulnerabilities.



Answering RQ1 Cont'd



Finding 2:

MulanPSL-2.0 has zero number of CVEs as a permissive license.
However, EPL-2.0, has high number of CVEs as a restrictive license.



Answering RQ1 Cont'd



Finding 3:

After removing libraries with both permissive and restrictive licenses for all licenses, EPL-2.0 went from **881 CVEs to 13 CVEs**.

Answering RQ1 Cont'd

License Name	Total Releases	License Name	Release Number
apache-2.0	30,822	apache-2.0	30,822
bsd-3-clause	4,276	bsd-3-clause	4,276
epl-1.0	8,743	epl-1.0	8,628
epl-2.0	6,143	epl-2.0	1,546
gpl-3.0	6,893	gpl-3.0	5,981
lgpl-2.1+	5,405	lgpl-2.1+	5,311
lgpl-3.0	3,764	lgpl-3.0	3,198
mit	15,011	mit	15,011
mulanpsl-2.0	4,050	mulanpsl-2.0	4,050

Finding 3:

After removing libraries with both permissive and restrictive licenses for all licenses, EPL-2.0 went from 881 CVEs to 13 CVEs.



WATERLOO ATTERNATION

Answering RQ2:

RQ2: Do specific license types **correlate with higher or lower** vulnerability incidence in Maven Central artifacts?



 Table 2. Updated number of releases after removing libraries

 with both permissive and restrictive licenses

License Name	Release Number
apache-2.0	30,822
bsd-3-clause	4,276
epl-1.0	8,628
epl-2.0	1,546
gpl-3.0	5,981
lgpl-2.1+	5,311
lgpl-3.0	3,198
mit	15,011
mulanpsl-2.0	4,050

Table 3. CVSS Severity Rating



Answering RQ2 Cont'd: Statistical test

Mann-Whitney U test:

- We choose the Mann-Whitney U test because it does not enforce any assumptions about the distribution of analyzed data and applies to a small number of data points.



Finding 4:

- The permissive licenses exhibit higher CVE incidence, and the Mann-Whitney U test further confirms that there is a statistically significant difference between permissive and restrictive licenses. Consequently, **our data provides** evidence of a correlation between license type and vulnerability incidence.



Threats to Validity

- Limited extraction (278K records out of 658K) might introduce selection bias.
- Analysis focused only on top 10 licenses and top 100 libraries per license, which may overlook less popular libraries.
- Temporal restrictions on CVE data (2020–2025) may not fully capture historical vulnerability trends.
- Weaver API's unidirectional query (from library to CVE) limits comprehensive vulnerability capture.

CS 846

PAGE 17



Future Work

- 1. Use full dataset to conduct analysis, whether through improved data extracting efficiency or devote more time into the process.
- 2. Expand analysis to include more license types and alternative ranking metrics.
- 3. Extend temporal range to capture longer CVE history.
- 4. Enhance API capabilities to support reverse querying (from CVE to libraries) for deeper insights.
- 5. Future studies could delve deeper into the implications of libraries that adopt multiple licenses, particularly those mixing permissive and restrictive types.



PRESENTATION TITLE

PAGE 19